

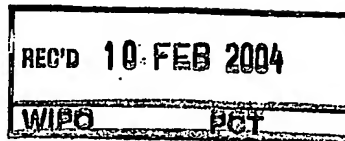


PCT/GB 2004 / 0 0 0 1 4 3



INVESTOR IN PEOPLE

**PRIORITY
DOCUMENT**
SUBMITTED OR TRANSMITTED IN
COMPLIANCE WITH RULE 17.1(a) OR (b)



The Patent Office
Concept House
Cardiff Road
Newport
South Wales
NP10 8QQ

I, the undersigned, being an officer duly authorised in accordance with Section 74(1) and (4) of the Deregulation & Contracting Out Act 1994, to sign and issue certificates on behalf of the Comptroller-General, hereby certify that annexed hereto is a true copy of the documents as originally filed in connection with the patent application identified therein.

In accordance with the Patents (Companies Re-registration) Rules 1982, if a company named in this certificate and any accompanying documents has re-registered under the Companies Act 1980 with the same name as that with which it was registered immediately before re-registration save for the substitution as, or inclusion as, the last part of the name of the words "public limited company" or their equivalents in Welsh, references to the name of the company in this certificate and any accompanying documents shall be treated as references to the name with which it is so re-registered.

In accordance with the rules, the words "public limited company" may be replaced by p.l.c., plc, P.L.C. or PLC.

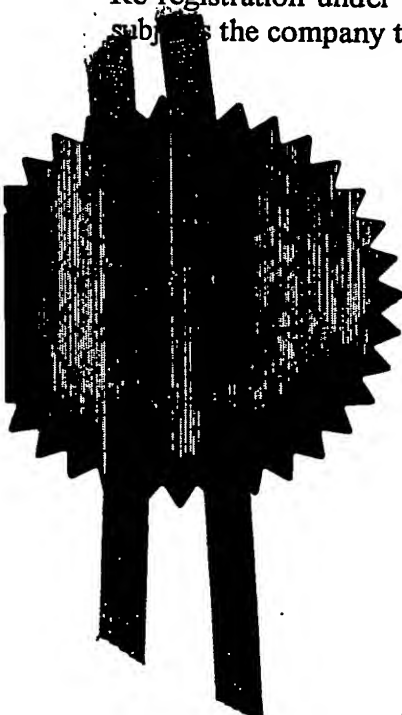
Re-registration under the Companies Act does not constitute a new legal entity but merely subjects the company to certain additional company law rules.

Signed

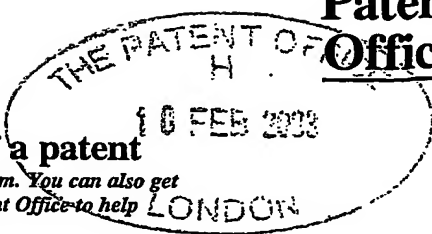
JEVENS

Dated 3 February 2004

BEST AVAILABLE COPY



Patents Act 1977
(Rule 16)



Request for grant of a patent

(See the notes on the back of this form. You can also get an explanatory leaflet from the Patent Office to help you fill in this form)

The Patent Office
Cardiff Road
Newport
Gwent NP10 8QQ

1.	Your reference	A30282	170 FEB 2003
2.	Patent application number (The Patent Office will fill in this part)	0303018.6	11 FEB 2003 170 FEB 2003 P01/7700 0.00-0303018.6
3.	Full name, address and postcode of the or of each applicant (<i>underline all surnames</i>)	BRITISH TELECOMMUNICATIONS public limited company 81 NEWGATE STREET LONDON, EC1A 7AJ, England Registered in England: 1800000	
	Patents ADP number (<i>if you know it</i>)	4867002	6300388001
	If the applicant is a corporate body, give the country/state of its incorporation	UNITED KINGDOM	
4.	Title of the invention	INFORMATION RETREIVAL	
5.	Name of your agent (<i>if you have one</i>)	NASH, Roger William	
	"Address for Service" in the United Kingdom to which all correspondence should be sent (<i>including the postcode</i>)	BT GROUP LEGAL SERVICES INTELLECTUAL PROPERTY DEPARTMENT HOLBORN CENTRE 120 HOLBORN LONDON, EC1N 2TE	
	Patents ADP number (<i>if you know it</i>)	4867001	7980311001
6.	If you are declaring priority from one or more earlier patent applications, give the country and the date of filing of the or of each of these earlier applications and (<i>if you know it</i>) the or each application number	Country	Priority application number (<i>if you know it</i>)
			Date of filing (<i>day / month / year</i>)
7.	If this application is divided or otherwise derived from an earlier UK application, give the number and the filing date of the earlier application	Number of earlier application	Date of filing (<i>day/month/year</i>)
8.	Is a statement of inventorship and of right to grant of a patent required in support of this request? (<i>Answer 'Yes' if:</i> a) any applicant named in part 3 is not an inventor, or b) there is an inventor who is not named as an applicant, or c) any named applicant is a corporate body. (See note (d))		
	YES		

Patents Form 1/77

Enter the number of sheets for any of the following items you are filing with this form.
Do not count copies of the same document

Continuation sheets of this form

Description 12 ✓

Claim(s) 5 ✓

Abstract 1 ✓

Drawing(s) 5 +5

10. If you are also filing any of the following, state how many against each item

Priority Documents

Translations of priority documents

Statement of inventorship and right to grant of a patent (*Patents Form 7/77*)

Request for preliminary examination and search (*Patents Form 9/77*) YES ✓

Request for substantive examination (*Patents Form 10/77*)

Any other documents (please specify)

11. I/We request the grant of a patent on the basis of this application.
Signature(s) Date:

10 February 2003

NASH, Roger William, Authorised Signatory

12. Name and daytime telephone number of person to contact in the United Kingdom

Samantha Radley

020 7492 8146

Warning

After an application for a patent has been filed, the Comptroller of the Patent Office will consider whether publication or communication of the invention should be prohibited or restricted under Section 22 of the Patents Act 1977. You will be informed if it is necessary to prohibit or restrict your invention in this way. Furthermore, if you live in the United Kingdom, Section 23 of the Patents Act 1977 stops you from applying for a patent abroad without first getting written permission from the Patent Office unless an application has been filed at least 6 weeks beforehand in the United Kingdom for a patent for the same invention and either no direction prohibiting publication or communication has been given, or any such direction has been revoked.

Notes

- a) If you need help to fill in this form or you have any questions, please contact the Patent Office on 0645 500505.
- b) Write your answers in capital letters using black ink or you may type them.
- c) If there is not enough space for all the relevant details on any part of this form, please continue on a separate sheet of paper and write "see continuation sheet" in the relevant part(s). Any continuation sheet should be attached to this form.
- d) If you have answered 'Yes' *Patents Form 7/77* will need to be filed.
- e) Once you have filled in the form you must remember to sign and date it.
- f) For details of the fee and ways to pay please contact the Patent Office.

INFORMATION RETREIVAL

This invention relates to information retrieval and in particular to a method and apparatus for identifying and retrieving information taking account of a level of expertise likely to be required of a user accessing it, and to a particular method and apparatus for determining the level of expertise applicable to a given set of information.

It is known to classify documents according to a number of different criteria, in particular according to information topic. Numerous prior art techniques have been devised to achieve automatic or semi-automatic classification of documents. Known classification techniques have been applied in particular to information retrieval arrangements to group or to help locate documents relating to particular topics of interest. However, while a search for relevant documents may be successful in locating a number of documents relevant to a particular topic of interest, the intended audience for each document will vary and many located documents may prove unsuitable for particular users, being for example too general for a specialised user having significant expertise in the topic.

According to a first aspect of the present invention there is provided a method for determining a measure of the level of expertise applicable to an information data set, comprising the steps of:

- (i) selecting, in respect of each of a plurality of predetermined levels of expertise, a representative sample set of information data sets;
- (ii) determining, for each of said selected information data sets, the value of a metric indicative of the incidence, in a reference corpus of information, of terms comprised in the selected data set; and
- (iii) using the values of said metric determined in step (ii) to train an information classifier to identify at least one of said plurality of predetermined levels of expertise applicable to an information data set using a value of said metric determined for the information data set.

The metric chosen for use in preferred embodiments of the present invention has the property that the values of the metric, calculated for different representative samples of data sets in a training set selected in step (i) above, fall within substantially distinct ranges. This enables a document classifier to be trained to rate

a given information data set according to which of the predetermined levels of expertise is most applicable, based solely upon the value of the metric calculated for the information data set being rated.

A value for the metric is calculated with reference to a reference corpus of information in a relevant language. In preferred embodiments of the present invention, the reference corpus used is the British National Corpus, referenced below, although an equivalent corpus may be available in respect of languages other than English. The reference corpus provides a measure, for each term, of the incidence of that term in the language represented by the corpus. For the purposes of the present patent application, "term" is intended to relate to a word or phrase or part of a word, e.g. a stemmed word. Different more specialised corpi of information may be selected, for example a corpus representative of the use of terms in speech, a corpus representative of written use, or a corpus of children's literature in a particular language.

Preferably the metric comprises a combined measure of the incidence within an information data set of terms comprised in the information data set and of the incidence of each said term in the reference corpus. In this way, the observed incidence of a particular term in the reference corpus may be weighted more highly, and hence contribute more to the value of the metric, the more frequently that term is found to occur in the information data set being rated. A preferred formula for calculating values for the metric is given in the detailed description below.

Preferably, training the classifier comprises:

- (a) making distributions of normalised values of said metric for data sets in each of the representative sample sets selected at step (i), above; and
- (b) for each of said predetermined levels of expertise, identifying from said distributions a corresponding range of normalised values of said metric.

Normalised values of the metric are obtained, in a preferred embodiment of the present invention, by taking account of the length of the information data set being rated in comparison with the mean length of data sets used to construct the reference corpus.

In a preferred embodiment of the present invention, the trained classifier is arranged to determine a measure of the probability that a particular one of said predetermined levels of expertise is applicable to the information data set being rated.

For example, if it is found that distributions of the calculated values of the metric for the training samples of data sets are overlapping to some degree, then there may be more than one level of expertise yielding a non-zero probability of association with information data set being rated. An output expressed in the form of probabilities for
5 each predetermined level of expertise may be particularly useful in fuzzy processing arrangements.

Preferably, determining a value for said metric comprises applying a stemming algorithm to stem terms comprised in a respective information data set and determining the incidence of the stemmed terms in the reference corpus. In
10 particular, a algorithm such as Porter, M.F., 1980, "An algorithm for suffix stripping", *Program*, 14(3) :130-137, since reprinted in Sparck Jones, Karen, and Peter Willet, 1997, *Readings in Information Retrieval*, San Francisco: Morgan Kaufmann, ISBN 1-55860-454-4, may be used to stem terms prior to obtaining their measure of incidence in the reference corpus.

15 According to a second aspect of the present invention there is provided a method of accessing information data sets, stored in an information system, relevant to search criteria specifying an indication of a category of information to be accessed and an indication of a predetermined level of expertise in respect of said category of information, the method comprising the steps of:

20 (i) selecting a training set of information data sets comprising, for each of a predetermined plurality of levels of expertise, a representative sample set of information data sets;

(ii) determining, for each data set in the training set, the value of a metric indicative of the incidence, in a reference corpus of information, of terms comprised
25 in the training data set;

(iii) using the values of said metric determined in step (ii) to train an information classifier to identify at least one of said predetermined plurality of levels of expertise applicable to a given information data set;

(iv) applying an information searching algorithm to identify information data
30 sets stored in said information system relevant to said specified category of information; and

(v) using the classifier trained at step (iii) to determine respective levels of expertise for information data sets identified at step (iv) and comparing the

determined levels of expertise with the level of expertise specified in said search criteria to thereby select relevant information data sets.

When searching for documents relevant to a particular category of information, by taking account also of the level of expertise of a user initiating the search in that information category and matching the user's level of expertise with that determined as being necessary for documents identified in the search, the search results selected for presentation to that user are likely to be more useful than those in a similar arrangement that otherwise ignores the intended level of expertise of readers of identified documents.

10 According to a third aspect of the present invention there is provided an apparatus for determining a level of expertise applicable to an information data set, the level of expertise being selected from a predetermined plurality of levels of expertise, the apparatus comprising:

an input for receiving an information data set;

15 calculating means arranged with access to a reference corpus of information to calculate, for an information data set, the value of a metric indicative of the incidence, in the reference corpus, of terms comprised in the information data set;

a trainable classifier; and

training means for training said classifier to identify, using a training set of information data sets comprising, for each of said predetermined plurality of levels of expertise, a representative sample set of information data sets and respective values of said metric, an applicable level of expertise selected from said predetermined plurality of levels of expertise for a received information data set;

20 wherein, in operation, on receipt of an information data set at said input, said calculating means are arranged to calculate a respective value for said metric and to input the calculated value to said trainable classifier, trained by said training means, to determine and output an indication of at least one of said predetermined plurality of levels of expertise applicable to said received information data set.

According to a fourth aspect of the present invention there is provided an information retrieval apparatus for accessing information data sets, stored in an information system, relevant to received search criteria specifying an indication of a category of information to be accessed and an indication of a predetermined level of expertise in respect of said category of information, the apparatus comprising:

30

calculating means arranged with access to a reference corpus of information to calculate, for an information data set, the value of a metric indicative of the incidence, in the reference corpus, of terms comprised in the information data set;

a trainable classifier;

5 training means for training said classifier to identify, using a training set of information data sets comprising, for each of a predetermined plurality of levels of expertise, a representative sample set of information data sets and respective values of said metric, an applicable level of expertise selected from said predetermined plurality of levels of expertise for a given information data set;

10 searching means for identifying information data sets in said information system relevant to said specified category of information to be accessed; and

selecting means arranged to trigger said calculating means to calculate values of said metric for information data sets identified by said searching means, to input the values so calculated to said trainable classifier, trained by said training
15 means, to determine and output respective applicable levels of expertise selected from said predetermined plurality of levels of expertise, and to select, for access, information data sets from those identified by said searching means having respectively determined levels of expertise that match said specified level of expertise.

20 According to a fifth aspect of the present invention there is provided an information retrieval apparatus for accessing information data sets, stored in an information system, relevant to received search criteria specifying an indication of a category of information to be accessed and to a specified indication of a predetermined level of expertise in respect of said category of information, the
25 apparatus comprising:

calculating means arranged with access to a reference corpus of information to calculate, for an information data set, the value of a metric indicative of the incidence, in the reference corpus, of terms comprised in the information data set;

an information classifier, trained, using, for each of a plurality of
30 predetermined levels of expertise, a representative sample set of training information data sets and respective values of said metric, to determine a level of expertise, selected from said plurality of predetermined levels of expertise, applicable to an information data set;

searching means for identifying information data sets in said information system relevant to said specified category of information to be accessed; and

selecting means arranged to trigger said calculating means to calculate values of said metric for information data sets identified by said searching means, to
5 input the values so calculated to said information classifier to determine and output respective applicable levels of expertise selected from said plurality of predetermined levels of expertise, and to select, for access, information data sets from those identified by said searching means having respectively determined levels of expertise that match said specified level of expertise.

10 A apparatus according to the fifth aspect of the present invention may be supplied with a ready-trained information classifier rather than one that has yet to be trained. An information classifier already trained using a general cross-section of training information data sets has been found to provide an acceptable level of performance when used to access information data sets across a range of
15 information categories.

Preferred embodiments of the present invention will now be described, by way of example only, with reference to the accompanying drawings of which:

Figure 1 is a diagram showing a trainable document classifier usable in an apparatus according to a first embodiment of the present invention;

20 Figure 2 is a diagram showing typical distributions of a preferred metric for a training sample of documents;

Figure 3 is flow diagram showing steps in a preferred training process;

Figure 4 is a flow diagram showing preferred steps in operation of the apparatus of Figure 1; and

25 Figure 5 is an information retrieval apparatus according to a second embodiment of the present invention.

This invention arises from the observation by the inventors in the present case that a metric comprising a statistical measure of the "commonality" of terms occurring in a document with reference to a corpus of information representative of
30 the use of words in a particular language can be used to train a conventional document classifier to distinguish those documents intended for general readership from those directed to a more expert reader. In the English language in particular, this metric may be calculated preferably with reference to the British National Corpus – a

100,000,000 word electronic databank sampled from the whole range of present-day English, spoken and written. Word frequencies for the British National Corpus have been published for example in "Word Frequencies in Written and Spoken English: based on the British National Corpus." by Geoffrey Leech, Paul Rayson and Andrew
 5 Wilson, published (2001) by Longman, London, ISBN 0582-32007-0 (Paperback).

A first embodiment of the present invention will now be described with reference to Figure 1.

Referring to the diagram of Figure 1, a trained document classifier 100 is shown that has been trained, by a process to be described below, to determine and
 10 to output a rating corresponding to one of a number of predefined levels of expertise to be associated with a given document 105, or to determine and to output a probability that the given document 105 relates to one or more of those predefined levels of expertise. A metric calculator 110 is arranged with access to a reference corpus 115 of information in a particular language to enable it to calculate, for the
 15 given document 105, the value of a metric, to be defined below, indicative of the "commonality" of terms occurring in the document 105. The classifier 100 has been trained to use a value of the metric calculated by the metric calculator 110 to determine the appropriate level of expertise to associate with the document 105. The expertise rating output by the trained classifier 100 may be used in a number of
 20 different applications, in particular in an improved information retrieval arrangement where only those documents that match a user's measure of expertise in a particular field of information are selected from a set of search results for presentation to the user.

A preferred metric found to be suitable for use with a document classifier
 25 100 to determine an expertise rating for a given document 105 is derived as follows.

A value α is first calculated, by the metric calculator 110, for the given document 105 using the formula

$$\alpha = \sum_i \log(tf_i + 1) \log\left(\frac{n(i)}{N}\right)$$

where tf_i is the term frequency within the given document 105 of the i -th distinct (preferably stemmed using the algorithm referenced above) term of the given document 105,

$n(i)$ is the number of documents in the reference corpus 115 containing the i -th distinct (stemmed) term of the given document 105 and

N is the total number of documents in the reference corpus 115.

Preferably the value of $n(i)/N$ is available directly as output from an interface to the reference corpus 115 for any particular stemmed term. For example, for a particular stemmed term, the reference corpus 115 returns a value representing the frequency with which the particular stemmed term occurs per million terms in the corpus 115.

The preferred metric then calculated by the metric calculator 110 is a "normalised" value for α , obtained by dividing α by a value β , where β is defined by:

$$\beta = \frac{\text{length_of_the_given_document}}{\text{mean_length_of_documents_in_the_reference_corpus}}$$

It has been found that when the values for this preferred metric α/β are plotted for a range of documents, those documents typically directed to "expert" readers in a particular field have a substantially distinct range of values for α/β in comparison with that for documents intended for more "general" readership. The differences in the two distributions can be seen, for a particular sample of documents, in Figure 2.

Referring to Figure 2, two distributions are shown, one distribution 200 for a sample of documents known to be intended for "general" readership and one distribution 205 for a sample of documents known to be intended for more "expert" readership. If more than two levels of expertise are to be distinguished, then samples of documents may be selected representative of one or more intermediate levels of expertise and the corresponding distributions plotted. Distributions may also be made in respect of samples of documents distinguishing "child" from "adult" levels of "expertise".

There are numerous variations to the formulae provided above for calculating α and β of the preferred metric, for use in preferred embodiments of the

present invention, that would be apparent to a person of ordinary skill, each variation taking account of the "commonality" of terms occurring within a given document. In addition, there are numerous variations in the way in which terms of a given document may be selected for use in calculating a value for the preferred metric. For
5 example, rather than considering every term within a given document, a known algorithm may be used to select terms most likely to be indicative of the information content of the given document, for example an algorithm to extract so-called "key terms" as described in European patent number EP 1032896 by the present Applicants. In a further variation, the reference corpus 115 used in preferred
10 embodiments of the present invention may be selected from a range of specialised corpi according to the particular information topic of documents under consideration or, more generally, according to whether the documents under consideration relate to technical or non-technical subject matter, or to children's literature for example.

Having determined a suitable metric as defined above, the next step is to use
15 that metric to train a document classifier either to identify which of the predefined levels of expertise to associate with a given document 105, or to determine a set of probabilities that a given document 105 is associated with one or more of the predefined levels of expertise. To this end, steps in a preferred training process will now be described with reference the flow diagram of Figure 3.

20 Referring to Figure 3, the training process begins with, at STEP 300, selection of a training set of documents comprising, for each of the predetermined levels of expertise to be applied, a representative training sample of documents known to contain subject matter expressed in a way suitable for readers having that level of expertise, e.g. "expert" readers or those with only a "general" appreciation of
25 a given information topic. In practice, while the training set of documents may relate to a particular information topic and a different training set of documents may be selected for each information topic, it has been found that a more general training set yields acceptable results when used to rate documents relating to a number of different information topics. At STEP 305, the value for the preferred metric α/β is
30 calculated, for example by the metric calculator 110, for each of the documents in the training set. At STEP 310, knowing the level of expertise associated with each document of the training set and the corresponding values for α/β , a conventional document classifier is trained to associate a given document 105 with one of the

predefined levels of expertise on the basis of a respective value for α/β . Preferably, the document classifier may be trained at STEP 310 by making distributions of document frequency in the respective training sample sets for values of α/β , as in Figure 2, and on the basis of the document frequency distributions for each sample, 5 determining the range of values of α/β corresponding to each of the pre-defined levels of expertise (there being two levels of expertise – "General" and "Expert" - in the example of Figure 2). Alternatively, if required, the document classifier 100 may be arranged, after training, to output probability values in respect of each of the predefined levels of expertise yielding a non-zero probability for the given document 10 105.

Steps in a preferred process, operable by the apparatus of Figure 1, for determining the level of expertise for a given document 105, will now be described with reference to the flow diagram of Figure 4.

Referring to Figure 4, the preferred process begins at STEP 400 with receipt 15 of a document 105 to be rated. At step 405 the value of the preferred metric α/β is calculated by the metric calculator 110 for the received document 105 using the formulae provided above, with reference to the reference corpus 115. Preferably, when accessing the reference corpus 115 to obtain a relative frequency score for a stemmed form of a particular term, if the reference corpus 115 provides relative 20 frequency scores for homonyms of the particular term, the metric calculator 110 is arranged to sum the relative frequencies provided for each homonym. That is, no attempt is made by the metric calculator 110 to distinguish use of a particular term in a given document 105 as a preposition from its use as an adjective, for example, before obtaining the relative frequency score from the reference corpus 115. 25 However, the metric calculator 110 may be arranged optionally to implement a known algorithm to analyse terms in the given document 105 and to identify the particular use of each term before obtaining the respective score for that use of the term from the reference corpus 115.

The resultant value for α/β is input, at STEP 410, to the trained document 30 classifier 100, preferably trained according to the process of Figure 3, and at STEP 415 the trained document classifier 100 outputs either an indication of the level of expertise to associate with the received document 105 or a set of probabilities that

the received document 105 is associated with each of one or more of the levels of expertise. This latter output is of particular use in fuzzy processing systems.

A preferred information retrieval apparatus will now be described with reference to Figure 5, incorporating the trained document classifier 100 of Figure 1 in
5 a preferred embodiment of the present invention.

Referring to Figure 5, an information retrieval software agent 500 is arranged to operate on behalf of a user to identify documents relevant to the user's submitted search criteria 505. Search criteria 505 typically comprise a set of keywords/phrases relating to a particular category of information sought by the user. The information
10 retrieval software agent 500 is arranged with access to a user profile store 510 wherein a predefined user profile may be stored for the user, the profile containing an indication of the level of expertise of the user in respect of the particular category of information being sought. However, the level of expertise of the user submitting the search criteria 505 may optionally be specified within the search criteria 505, so
15 obviating the need for the information retrieval software agent 500 to make a separate access to the user profile store 510 to obtain the user's expertise level.

The information retrieval software agent 500 is arranged with access to the Internet 515 and hence to one or more search engines 520 to help identify and retrieve sets of information stored on web servers 525 relevant to the user's
20 submitted search criteria 505. The information retrieval software agent 500 is also arranged with access to a trained document classifier 100 as above, by way of a metric calculator 110 arranged with access to a reference corpus 115 for calculating a value for the metric α/β , as defined above, for a particular document, which value when input to the trained document classifier 100 enables the level of expertise
25 associated with the particular document to be determined. The information retrieval software agent 500 is arranged to output a list of search results 530 in response to the user's submitted search criteria 505, the search results 530 being tailored both to the user's specified category of information (505) and to the user's level of expertise (510) with respect to that category of information (505).

30 In operation, the information retrieval software agent 500 is arranged, on receipt of search criteria 505 submitted by a user, to access the user's personal profile 510 to determine the level of expertise of the user in respect of the category of information represented by the submitted criteria 505, assuming that the user has

not specified his/her level of expertise within the search criteria 505. The information retrieval software agent 500 then accesses search engines 520 or web servers 525 directly to identify and retrieve sets of information relevant to the information category specified in the submitted search criteria 505, by conventional means. As
5 relevant information sets are identified and received, the information retrieval software agent 500 determines the level of expertise to be associated with each relevant information set using functionality provided by the metric calculator 110 and the trained document classifier 100, as described above with reference to Figure 4. The information retrieval software agent 500 compares the level of expertise
10 determined for each relevant information set with the level of expertise (510) of the user and thereby selects, to output to the user as search results 530, a set of relevant information sets having determined levels of expertise matching the user's level of expertise.

In a further embodiment of the present invention a trained document
15 classifier 100 may be used to derive a measure of the level of expertise of a user in respect of a particular information topic. By monitoring information retrieval activity of a user in respect of the particular information topic, those documents that the user evidently finds useful, for example because the user retrieves a whole document to read or provides feedback as to the usefulness of the document, may be input to the
20 metric calculator 110 and the respective metric values input to the trained document classifier 100 to determine the level of expertise to associate with these "useful" documents and hence, by implication, the level of expertise of the user in the information topic that those documents represent.

It would be apparent to a person of ordinary skill in this field of information
25 retrieval, that preferred embodiments of the present invention may be applied in other information retrieval arrangements in which the expertise of a user may be taken into account when selecting information for presentation to that user or otherwise used in respect of that user.

CLAIMS

1. A method for determining a measure of the level of expertise applicable to a given information data set, comprising the steps of:
 - 5 (i) selecting, in respect of each of a plurality of predetermined levels of expertise, a representative sample set of information data sets;
 - (ii) determining, for each of said selected information data sets, the value of a metric indicative of the incidence, in a reference corpus of information, of terms comprised in the selected information data set; and
 - 10 (iii) using the values of said metric determined in step (ii) to train an information classifier to identify, from a value of said metric calculated for the given information data set, at least one of said plurality of predetermined levels of expertise applicable to the given information data set.
- 15 2. A method as in Claim 1, wherein said metric comprises a combined measure of the incidence within an information data set of terms comprised in the information data set and of the incidence of each said term in the reference corpus.
3. A method as in Claim 1 or Claim 2, wherein at step (iii), training the classifier
20 comprises:
 - (a) making distributions of normalised values of said metric for data sets in each of the representative sample sets selected at step (i); and
 - (b) for each of said predetermined levels of expertise, identifying from said distributions a corresponding range of normalised values of said metric.
- 25 4. A method as in any one of claims 1 to 3, wherein at step (iii), the trained classifier is arranged to determine a measure of the probability that a particular one of said predetermined levels of expertise is applicable to the information data set.
- 30 5. A method as in any one of the preceding claims, wherein determining a value for said metric comprises applying a stemming algorithm to stem terms comprised in a respective information data set and determining the incidence of the stemmed terms in the reference corpus.

6. A method as in any one of the preceding claims, wherein the reference corpus is provided with an interface for outputting the relative frequency of occurrence in the corpus of a term.

5

7. A method of accessing information data sets, stored in an information system, relevant to search criteria specifying an indication of a category of information to be accessed and to a specified indication of a predetermined level of expertise in respect of said category of information, the method comprising the steps

10 of:

(i) selecting a training set of information data sets comprising, for each of a plurality of predetermined levels of expertise, a representative sample set of information data sets;

15 (ii) determining, for each data set in the training set, the value of a metric indicative of the incidence, in a reference corpus of information, of terms comprised in the training data set;

(iii) using the values of said metric determined in step (ii) to train an information classifier to identify at least one of said plurality of predetermined levels of expertise applicable to a given information data set;

20 (iv) applying an information searching algorithm to identify information data sets stored in said information system relevant to said specified category of information; and

25 (v) using the classifier trained at step (iii) to determine respective levels of expertise for information data sets identified at step (iv) and comparing the determined levels of expertise with the specified level of expertise to thereby select relevant information data sets.

8. An apparatus for determining a level of expertise applicable to an information data set, the level of expertise being selected from a plurality of predetermined levels of expertise, the apparatus comprising:

an input for receiving an information data set;

calculating means arranged with access to a reference corpus of information to calculate, for an information data set, the value of a metric indicative of the incidence, in the reference corpus, of terms comprised in the information data set;

a trainable classifier; and

5 training means for training said classifier to identify, using a training set of information data sets comprising, for each of said plurality of predetermined levels of expertise, a representative sample set of information data sets and respective values of said metric, an applicable level of expertise selected from said plurality of predetermined levels of expertise for a received information data set;

10 wherein, in operation, on receipt of an information data set at said input, said calculating means are arranged to calculate a respective value for said metric and to input the calculated value to said trainable classifier, trained by said training means, to determine and output an indication of at least one of said plurality of predetermined levels of expertise applicable to said received information data set.

15

9. An apparatus as in Claim 8, wherein said metric comprises a combined measure of the incidence within an information data set of terms comprised in the information data set and of the incidence of each said term in the reference corpus.

20 10. An apparatus as in Claim 8 or Claim 9, wherein said training means are arranged to train said trainable classifier using the steps of:

(a) making distributions of normalised values of said metric for data sets in each of the representative sample sets; and

(b) for each of said predetermined levels of expertise, identifying from said
25 distributions a corresponding range of normalised values of said metric.

11. An apparatus as in any one of claims 8 to 10, wherein said trainable classifier is arranged, after training by said training means, to determine a measure of the probability that a particular one of said plurality of predetermined levels of
30 expertise is applicable to a received information data set.

12. An apparatus as in any one of claims 8 to 11, wherein said calculating means are arranged to calculate a value for said metric by applying a stemming

algorithm to stem terms of a respective information data set and by determining the relative incidence of the stemmed terms in the reference corpus.

13. An information retrieval apparatus for accessing information data sets, stored
5 in an information system, relevant to received search criteria specifying an indication of a category of information to be accessed and to a specified indication of a predetermined level of expertise in respect of said category of information, the apparatus comprising:

calculating means arranged with access to a reference corpus of information
10 to calculate, for an information data set, the value of a metric indicative of the incidence, in the reference corpus, of terms comprised in the information data set;

a trainable classifier;

training means for training said classifier to identify, using a training set of
information data sets comprising, for each of a plurality of predetermined levels of
15 expertise, a representative sample set of information data sets and respective values of said metric, an applicable level of expertise selected from said plurality of predetermined levels of expertise for a given information data set;

searching means for identifying information data sets in said information
system relevant to said specified category of information to be accessed; and

20 selecting means arranged to trigger said calculating means to calculate values of said metric for information data sets identified by said searching means, to input the values so calculated to said trainable classifier, trained by said training means, to determine and output respective applicable levels of expertise selected from said plurality of predetermined levels of expertise, and to select, for access,
25 information data sets from those identified by said searching means having respectively determined levels of expertise that match said specified level of expertise.

14. An information retrieval apparatus for accessing information data sets, stored
30 in an information system, relevant to received search criteria specifying an indication of a category of information to be accessed and to a specified indication of a predetermined level of expertise in respect of said category of information, the apparatus comprising:

calculating means arranged with access to a reference corpus of information to calculate, for an information data set, the value of a metric indicative of the incidence, in the reference corpus, of terms comprised in the information data set;

an information classifier, trained, using, for each of a plurality of
5 predetermined levels of expertise, a representative sample set of training information data sets and respective values of said metric, to determine a level of expertise, selected from said plurality of predetermined levels of expertise, applicable to an information data set;

searching means for identifying information data sets in said information
10 system relevant to said specified category of information to be accessed; and

selecting means arranged to trigger said calculating means to calculate values of said metric for information data sets identified by said searching means, to input the values so calculated to said information classifier to determine and output
15 respective applicable levels of expertise selected from said plurality of predetermined levels of expertise, and to select, for access, information data sets from those identified by said searching means having respectively determined levels of expertise that match said specified level of expertise.

ABSTRACT

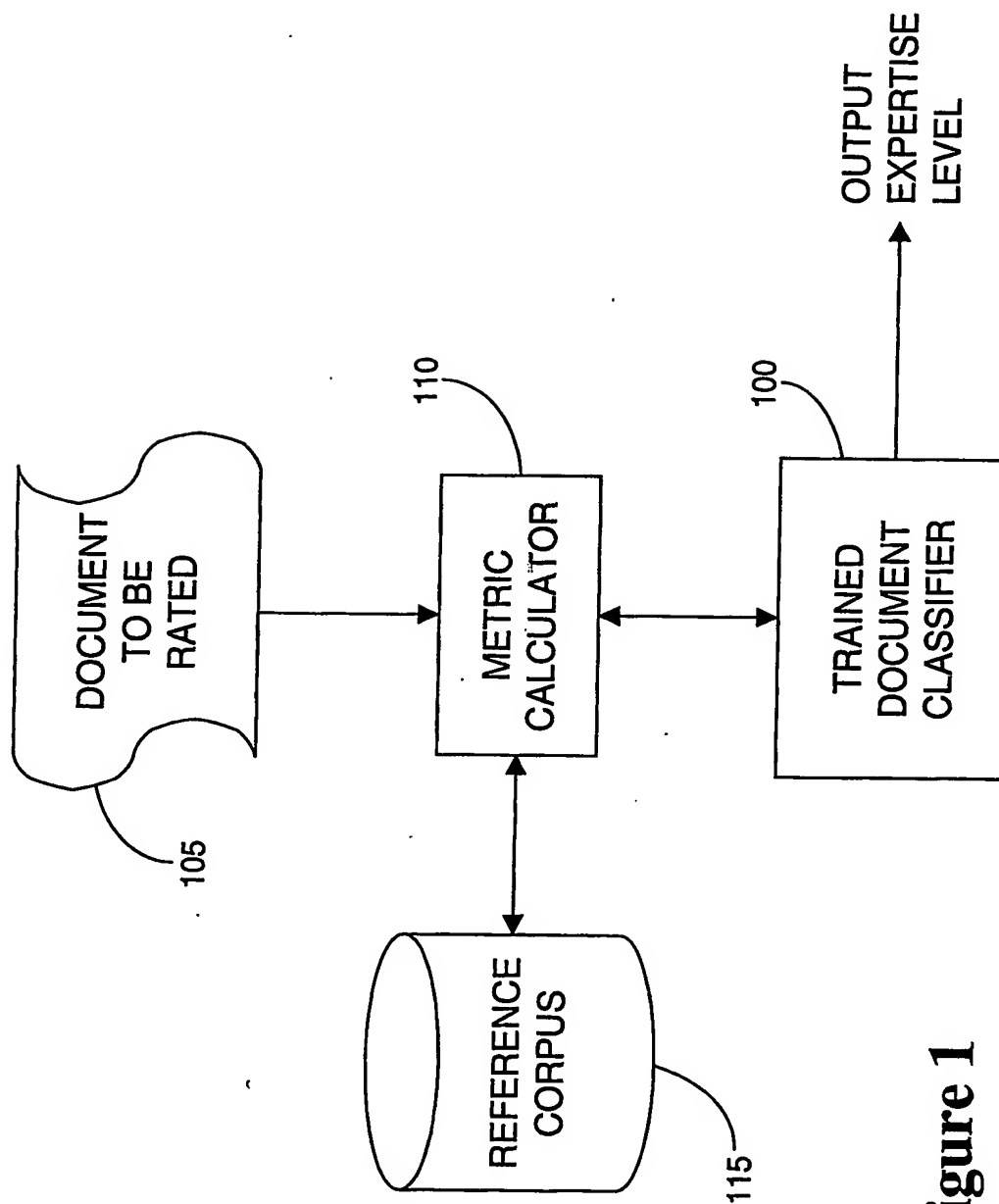
INFORMATION RETRIEVAL

5

An apparatus and method are provided for determining a level of expertise applicable to a particular document and for using this determined level of expertise in an improved information retrieval arrangement. A trainable document classifier is used to identify an applicable level of expertise using a metric indicative of the commonality, as measured with reference to a reference corpus, of terms comprised in a given document, trained using a training set of documents comprising, for each of a plurality of predetermined levels of expertise, a representative sample of documents and their respective metric values. An information retrieval apparatus is arranged to identify documents relevant to a specified category of information and to select from documents so identified those having a level of expertise, determined by the trained document classifier, matching a specified level of expertise for a target user in respect of that category of information.

Figure (1)

20

**Figure 1**

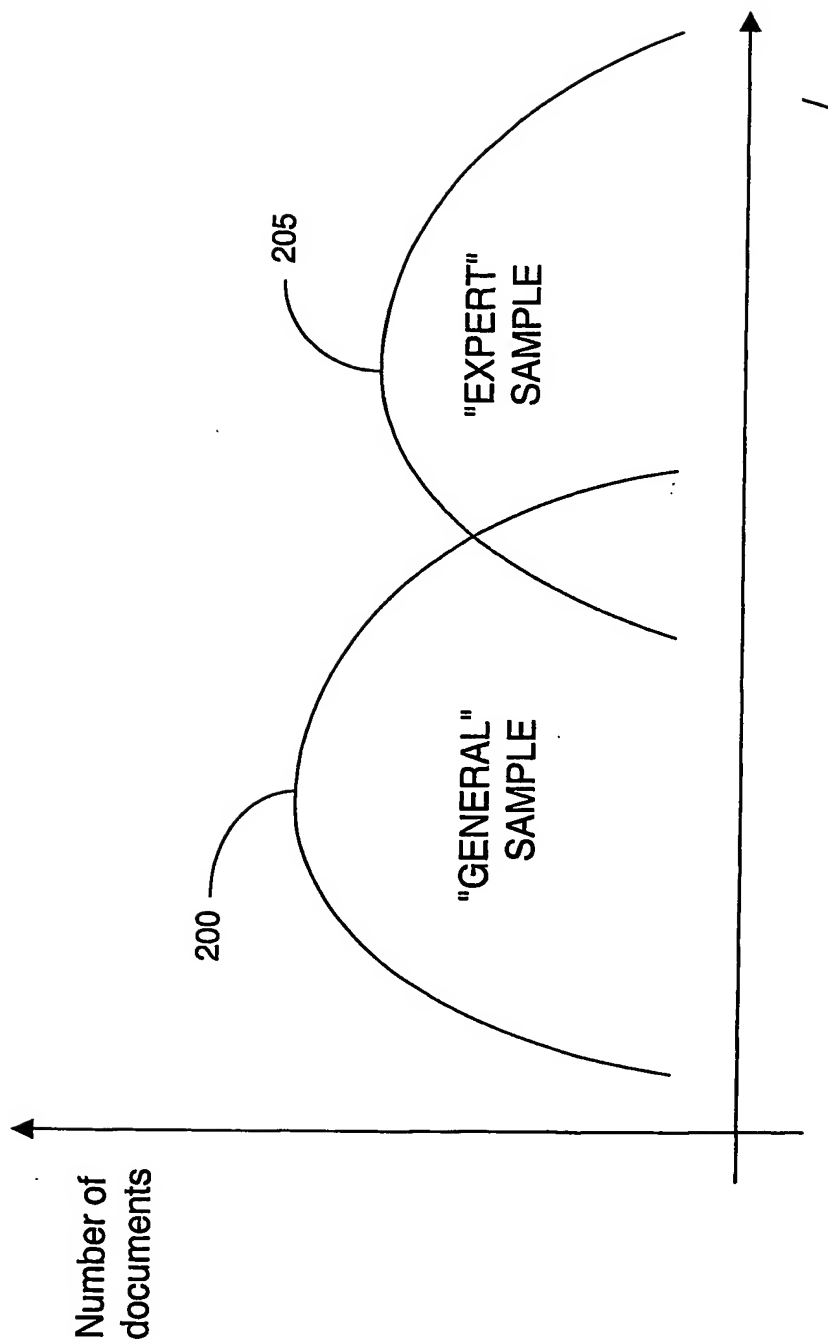
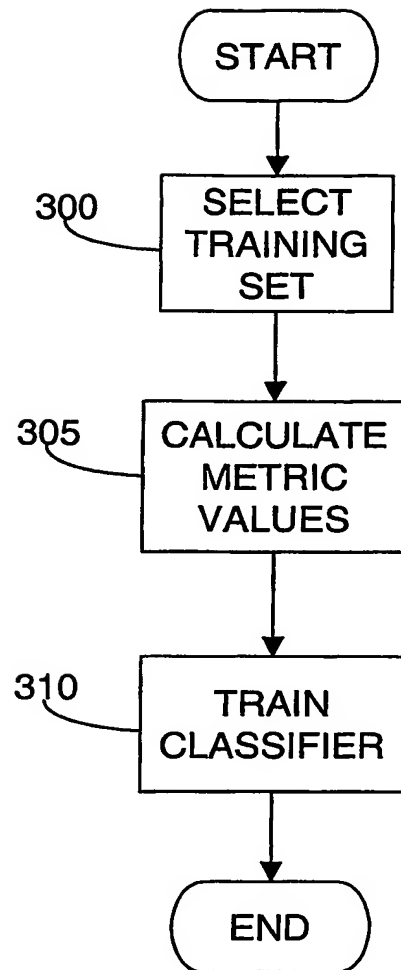
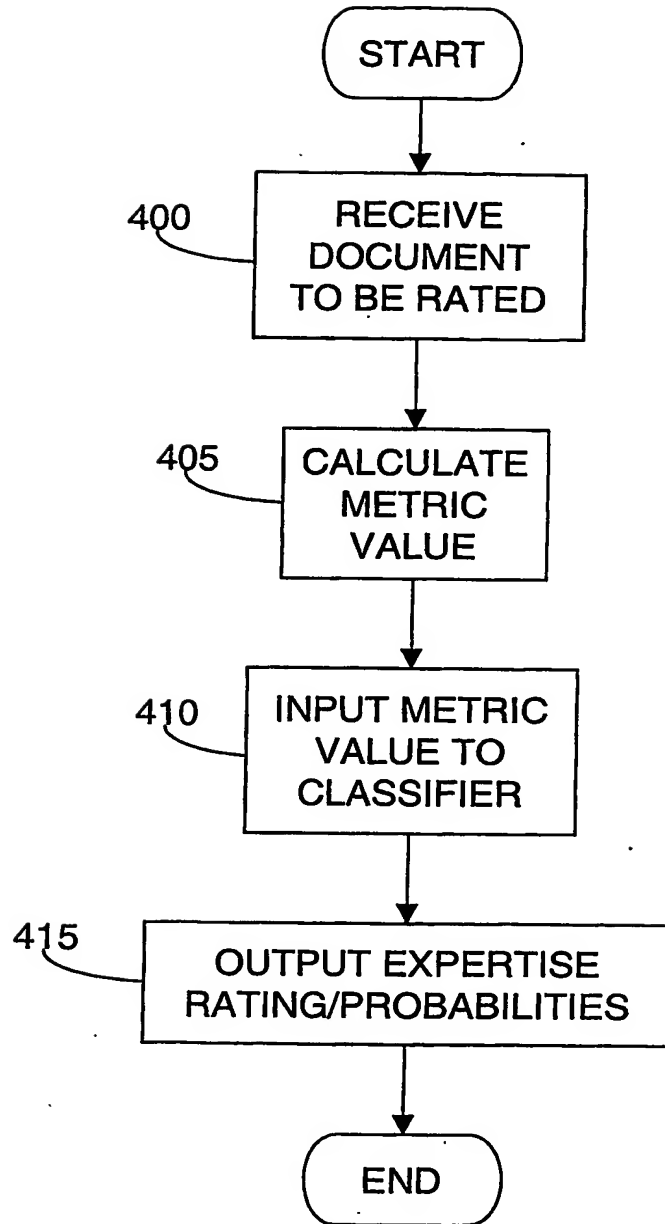
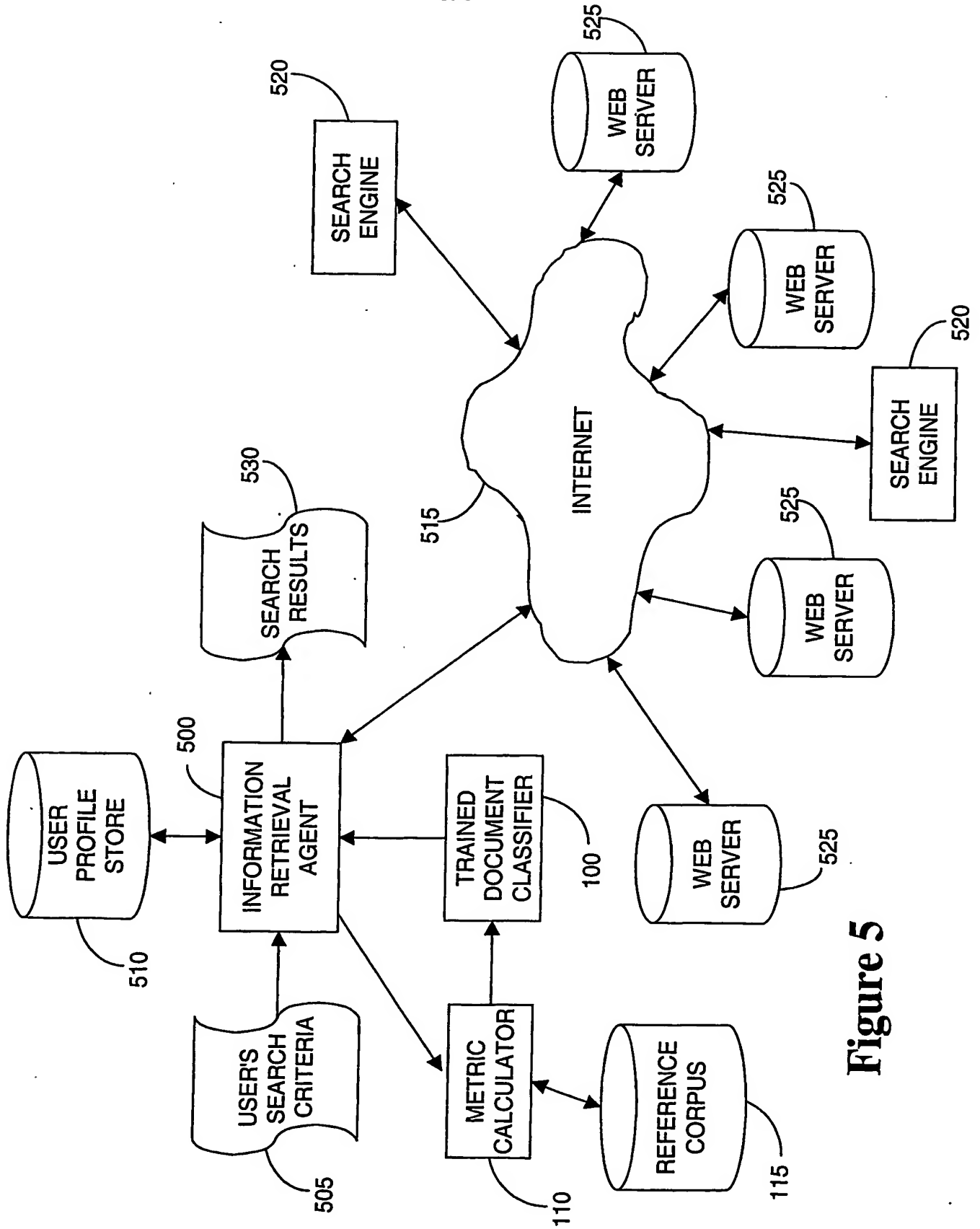


Figure 2

**Figure 3**

**Figure 4**

**Figure 5**

PCT Application
PCT/GB2004/000143



**This Page is Inserted by IFW Indexing and Scanning
Operations and is not part of the Official Record**

BEST AVAILABLE IMAGES

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images include but are not limited to the items checked:

- ☒ BLACK BORDERS
- ☒ IMAGE CUT OFF AT TOP, BOTTOM OR SIDES
- ☒ FADED TEXT OR DRAWING
- ☐ BLURRED OR ILLEGIBLE TEXT OR DRAWING
- ☐ SKEWED/SLANTED IMAGES
- ☒ COLOR OR BLACK AND WHITE PHOTOGRAPHS
- ☐ GRAY SCALE DOCUMENTS
- ☐ LINES OR MARKS ON ORIGINAL DOCUMENT
- ☐ REFERENCE(S) OR EXHIBIT(S) SUBMITTED ARE POOR QUALITY
- ☐ OTHER: _____

IMAGES ARE BEST AVAILABLE COPY.

As rescanning these documents will not correct the image problems checked, please do not report these problems to the IFW Image Problem Mailbox.